

Volker Roth

Rechnersicherheit, SoSe 21

Übung 03

TutorIn: Oliver Wiese

Tutorium 02

Materialien: Latex, VSC, Skript

26. Mai 2021

1 Shoulder surfer meets Markov

(a) Give an algorithm with observed digits (less than 4) as input and outputs three likely PINs.

1. We assume that we always observe only the first digit of the PIN.
2. First we preprocessed the Markov model that we expect to get thought a file. We expect this form of representation, because we defined it like that in the last exercise sheet. The following snippet shows the form we expect of a Markov model:

```

1  "00": 0.09766187759656811,
2  "01": 0.10964242399751109,
3  "02": 0.09521532801048535,
4  "03": 0.08118016683277238,
5  "04": 0.07375969673496395,
6  "05": 0.08015285944438737,
7  "06": 0.0798709585070059,
8  "07": 0.07797165702734245,
9  "08": 0.09059390538937141,
10 "09": 0.08103410417091668,
11 "0\u22a4": 0.11190444775411616,
12 "10": 0.13491879837536758,

```

Here is the Markov model after we preprocessed it. Now it is ordered from highest possibility that that char is taken to lowest:

```

1  "0": [
2      [
3          0.10964242399751109,
4          "1"
5      ],
6      [
7          0.09766187759656811,
8          "0"
9      ],
10     [
11         0.09521532801048535,
12         "2"

```

And here is the code on how we did it:

```
1 #returns a better dict format of markov_model
2 #following will be returned:
3 #eg: '0' : [('5',p1),('2',p2),('9',p3)...]
4 #    '1' : [('2',p1),('6',p2),('0',p3)...]
5 # where p1>p2>p3. so the possibilitys are sorted.
6 #that will allow us to get the 'best' edge if we search for
7 # an edge for a specific char (eg: result['1'][0][0] will return '2')
8 def preprocessing(markov_model):
9     #edge from v->e
10    #note that the edges of list_markov are sorted.
11    #eg '00':p, '01':p, '02':p, '03':p
12    # where p stands for a possibility
13    list_markov = list(markov_model)
14    v = list_markov[0][0]
15    temp = []
16    result = {}
17    for edge in list_markov:
18        #if a new different v1 (char) occurs
19        #we can save our current temp in the result
20        #and start to calc for the new v1(char)
21        #because the previous v will not occur again
22        #as list_markov edges are sorted
23        if v != edge[0]:
24            result[v] = sorted(temp, reverse=True)
25            v = edge[0]
26            temp = []
27
28        #append the possibility and e
29        #but only if e != top
30        #as our pin has a fixed len, we dont the end symbol
31        if edge[1:] != top:
32            temp.append( ( markov_model[edge], edge[1:] ) )
33
34    #ad the last digit (9) will be skipped otherwise
35    result[v] = sorted(temp, reverse=True)
36    return result
```

- After that we ran the following function. It takes the Markov model, the processed Markov model, the observed digit, the number of PINs we want to get and an accuracy int as input. The function is greedy and always takes the next digit with the highest possibility, which is quite easy, because of the preprocessing we did with the Markov model. After it found more PINs than 'num_of_output_pins' it drops the PIN with the lowest total possibility.

```

1 def guess_password(ordered_markov_model, markov_model ,accuracy, observed_digit,
2   num_of_output_pins):
3     result = []
4     for i_1 in range(accuracy):
5         #add the first (second) char to the pin and calc the probability
6         temp_prob_1, temp_char_1 = ordered_markov_model[observed_digit][i_1]
7         pin_1 = observed_digit + temp_char_1
8         prob_1 = 1 * temp_prob_1
9
10        for i_2 in range(accuracy):
11            #add
12            temp_prob_2, temp_char_2 = ordered_markov_model[temp_char_1][i_2]
13            pin_2 = pin_1 + temp_char_2
14            prob_2 = prob_1 * temp_prob_2
15
16            for i_3 in range(accuracy):
17                #add
18                temp_prob_3, temp_char_3 = ordered_markov_model[temp_char_2][i_3]
19                pin_3 = pin_2 + temp_char_3
20                prob_3 = prob_2 * temp_prob_3
21
22                #last char (digit) of the pin has to be the top symbol top = '\u22A4'
23                #so add that to the pin and calculate the probability of that edge
24                with the normal markov_model
25                pin_3 += top
26                prob_3 *= markov_model[temp_char_3+top]
27
28                #append result and drop lowest
29                result.append( (prob_3, pin_3) )
30                drop_lowest(result, num_of_output_pins)
31
32    return result

```

(b) Discuss the efficiency and success rate of your algorithm.

- The efficiency and success rate of your algorithm depends on the accuracy variable. As our algorithm has three for loops, its efficiency is accuracy^3 .
- Accuracy should/ can not be higher than the length of the alphabet, that is 10. So if we choose $\text{accuracy} = 10$ we basically try each possible combination.
- However with $\text{accuracy} = 10$ our algorithm always succeeds and returns the most likely 3 PINs. When we lower the accuracy we might not get the best 3 PINs, but in our tests half of that (5) seemed sufficient. What follows is a picture of the 3 most likely PINs of the Markov model from the RockYou dataset and how different accuracy values affect these PINs.

```

>>> main(2)
[(0.003192080901017929, '1234T'), (0.0008987235750735024, '1230T'), (0.0008493872205011611, '1212T')]
>>> main(4)
[(0.003192080901017929, '1234T'), (0.0009438533961205413, '1223T'), (0.0013149627524600077, '1123T')]
>>> main(6)
[(0.003192080901017929, '1234T'), (0.0025328192843752473, '1256T'), (0.002250275131115612, '1456T')]
>>> main(8)
[(0.003192080901017929, '1234T'), (0.0025328192843752464, '1256T'), (0.0022502751311156075, '1456T')]
>>> main(10)
[(0.003192080901017929, '1234T'), (0.002532819284375247, '1256T'), (0.0022502751311156066, '1456T')]

```

(c) Revisit the first-order Markov model from the previous assignment. What are your guesses if the adversary observes only a one.

- The expected input Markov model (snipped):

```

1  "00": 0,
2  "01": 0.16666666666666666,
3  "02": 0,
4  "03": 0.5,
5  "0\u22A4": 0.3333333333333333,
6  "10": 0.1111111111111111,
7  "11": 0,
8  "12": 0.2222222222222222,
9  "13": 0.4444444444444444,
10 "1\u22A4": 0.2222222222222222,
11 "20": 0.2222222222222222,
12 "21": 0,
13 "22": 0,

```

- After the preprocessing (snipped):

```

1  "0": [
2  [
3      0.5,
4      "3"
5  ],
6  [
7      0.16666666666666666,
8      "1"
9  ],
10 [
11    0,
12    "\u22A5"
13 ],
14 [
15    0,
16    "2"
17 ],
18 [
19    0,
20    "0"
21 ],
22 ],
23 "1": [
24 [
25     0.4444444444444444,
26     "3"
27 ],
28 [
29     0.2222222222222222,
30     "2"
31 ],
32 [
33     0.1111111111111111,
34     "0"
35 ],
36 [

```

- The output:

```

>>> main(2)
[(0.009259259259259, '1323T'), (0.006172839506172839, '1320T'), (0.006172839506172839, '1233T')]
>>> main(3)
[(0.009259259259259, '1323T'), (0.008333333333333333, '1303T'), (0.007716049382716049, '1230T')]
>>> main(4)
[(0.009259259259259, '1323T'), (0.008333333333333333, '1303T'), (0.007716049382716049, '1230T')]
>>> []

```

2 Shoulder surfer meets Markov II (Bonus)

(a) Create a markov model of 4- and 6-digit PINs based on the RockYoudata set using Python. You can either filter them by yourself or use the dataset of Markert et al.

- After some preprocessing of the RockYou dataset and error filtering (like in the last exercise sheet) we just filtered the dataset with a regex. Snipped of the code:

```
1 #create passwordset P1
2 def filter_4digits_6digits(df):
3     #\A from beginning of the string till \z end: looks for digits with 4-4 chars
4     #or 6-6 chars
5     a = '(?:\A(?:\d{4}\Z)|\A(?:\d{6}\Z))'
6     df = df.loc[df['passwords'].str.contains(a, regex=True)].copy()
7     #reset index
8     df.reset_index(inplace=True, drop=True)
9     return df
10
11 def main():
12     df = preprocess()
13     df = correct_errors(df)
14     df = filter_4digits_6digits(df)
15     write_to_file(df, output_file)
```

Snippet of the filtered dataset:

```
1 frequentcys passwords
2 0 290729 123456
3 1 13984 654321
4 2 13272 111111
5 3 13028 000000
6 4 9516 123123
7 5 7419 666666
8 6 5377 121212
9 7 5225 112233
10 8 5175 555555
11 9 5055 789456
12 10 5052 999999
13 11 4576 159753
14 12 4486 222222
15 13 4469 987654
16 14 4118 123321
17 15 3649 888888
18 16 3466 456789
19 17 3141 123654
20 18 3100 777777
```

(b) Implement your above algorithm in Python.

- First we load the dataset and add the bot and top symbol to every PIN

```
1 def init_database():
2     #read datafile
3     df = pd.read_csv(datafile, sep='\t', index_col = 0, dtype={'passwords': str})
4     #add top and bottom
5     df['passwords'] = bot + df['passwords'] + top
6     return df
```

- Then we count the occurrences of every char in the alphabet:

```

1 #count occurrences of chars in the password dataset
2 def count_occurrences(df, alphabet):
3     occurrences = {}
4     #for every char in alphabet (digit)
5     for char in alphabet:
6         #count the occurrence of that char, for every password in the database (*
        #by frequentcys)
7         df['occurrences'] = count_char(df['passwords'] , char) * df['frequentcys'
8     ]
9     #save result to dict
10    occurrences[char] = df['occurrences'].sum()
11    return occurrences
12
13 #count occurrence of a char in a string
14 def count_char(list, char):
15     result = []
16     for string in list:
17         result.append(string.count(char))
18     return result

```

```

occ: {'1': 2299585, '0': 2053913, '1': 2333574, '2': 1869088, '3': 1208374, '4': 1050749, '5': 1160702, '6': 1093210, '7': 77091
2, '8': 1120929, '9': 1094735, 'T': 2299585}

```

- After that we create the alphabet for the Markov model (the edges) and count the occurrences of these chars (edges) as well:

```

1 #returns an alphabet with every possible combination of the 2 input alphabets
2 # if you combine the chars in them
3 #the if clause will prevent that we get edges which will never occur
4 #e.g. T0, T1, ... TBot, ... OBot, 1Bot...
5 def markov_alphabet(alp1, alp2):
6     return [char1+char2 for char1 in alp1 for char2 in alp2 if char1 != top and
7             char2 != bot]

```

```

alphabet_1: ['1', '0', '1', '2', '3', '4', '5', '6', '7', '8', '9', 'T']

alphabet_2: ['10', '11', '12', '13', '14', '15', '16', '17', '18', '19', '1T', '00', '01', '02', '03', '04', '05', '06', '07', '
08', '09', '0T', '10', '11', '12', '13', '14', '15', '16', '17', '18', '19', '1T', '20', '21', '22', '23', '24', '25', '26', '27',
'28', '29', '2T', '30', '31', '32', '33', '34', '35', '36', '37', '38', '39', '3T', '40', '41', '42', '43', '44', '45', '46', '
47', '48', '49', '4T', '50', '51', '52', '53', '54', '55', '56', '57', '58', '59', '5T', '60', '61', '62', '63', '64', '65', '66',
'67', '68', '69', '6T', '70', '71', '72', '73', '74', '75', '76', '77', '78', '79', '7T', '80', '81', '82', '83', '84', '85', '
86', '87', '88', '89', '8T', '90', '91', '92', '93', '94', '95', '96', '97', '98', '99', '9T']

occ_edge_set: {'10': 495145, '11': 901009, '12': 330384, '13': 102202, '14': 69034, '15': 70438, '16': 75697, '17': 70500, '18':
102983, '19': 82193, '1T': 0, '00': 200589, '01': 225196, '02': 195564, '03': 166737, '04': 151496, '05': 164627, '06': 164048,
'07': 160147, '08': 186072, '09': 166437, '0T': 229842, '10': 314843, '11': 270394, '12': 629001, '13': 112622, '14': 112753, '15
': 105421, '16': 84293, '17': 89164, '18': 127976, '19': 206565, '1T': 217209, '20': 200787, '21': 210304, '22': 155452, '23': 45
6015, '24': 101217, '25': 127018, '26': 87732, '27': 85505, '28': 123614, '29': 104409, '2T': 194204, '30': 147571, '31': 127976,
'32': 97921, '33': 62935, '34': 334488, '35': 41146, '36': 44623, '37': 33771, '38': 51771, '39': 46387, '3T': 208526, '40': 964
01, '41': 82507, '42': 63799, '43': 53213, '44': 37362, '45': 353907, '46': 35807, '47': 38619, '48': 53536, '49': 43370, '4T': 1
84253, '50': 101754, '51': 87041, '52': 86141, '53': 49625, '54': 58953, '55': 53782, '56': 351858, '57': 37875, '58': 65032, '59
': 58922, '5T': 193496, '60': 92870, '61': 75007, '62': 56561, '63': 38717, '64': 27692, '65': 57512, '66': 58760, '67': 40000, '
68': 60099, '69': 66159, '6T': 498650, '70': 91789, '71': 73326, '72': 55537, '73': 31838, '74': 33840, '75': 39432, '76': 37963,
'77': 54659, '78': 91420, '79': 61854, '7T': 187976, '80': 116225, '81': 95484, '82': 79236, '83': 56231, '84': 59428, '85': 775
01, '86': 74460, '87': 95218, '88': 128198, '89': 132981, '8T': 189670, '90': 152781, '91': 121997, '92': 96661, '93': 66980, '94
': 56511, '95': 53695, '96': 56786, '97': 54176, '98': 113931, '99': 108707, '9T': 195759}

```

- And at last we calculate the Markov model with these occurrences:

```

1 #creates a makrov model / calcs the probability of the edge_set
2 def create_makrov_model(occ, alphabet, occ_edge_set):
3     markov_model = {}
4     for char in alphabet:
5         #divides the total num of edges by the total number of the beginning
        #vertice
6         #e.g. occ_edge_set['01'] = 10 / occ['0'] = 20 => 0.5
7         markov_model[char] = occ_edge_set[char] / occ[char[0]]
8     return markov_model

```

(c) *Evaluate your algorithms. What is the advantage of your algorithms compared to a random guess?*

- If we use the Markov model like in task 1. and output three likely PINs when we observe one digit. We get the following results:

```
>>> main(2)
[(0.003192080901017929, '1234T'), (0.0008987235750735024, '1230T'), (0.0008493872205011611, '1212T')]
>>> main(4)
[(0.003192080901017929, '1234T'), (0.0009438533961205413, '1223T'), (0.0013149627524600069, '1123T')]
>>> main(6)
[(0.003192080901017929, '1234T'), (0.0025328192843752464, '1256T'), (0.002250275131115608, '1456T')]
>>> main(8)
[(0.003192080901017929, '1234T'), (0.0025328192843752464, '1256T'), (0.002250275131115608, '1456T')]
>>> main(10)
[(0.003192080901017929, '1234T'), (0.0025328192843752464, '1256T'), (0.002250275131115608, '1456T')]
```

These PINs are better than a random guess, as they are much more likely to occur. A random guess would have the probability of $\frac{1}{10^3} = 0.001$, which is just a third compared to our most likely PIN: 0.0032.

3 Follow up

Please read Sections 4.1 and 4.2 of the paper³ and answer the following questions:

- (a) *What are the differences to the algorithm from the lecture?*
- (b) *Why do you think Algorithm 1 is different?*
- (c) *What are the consequences for storage and runtime?*

Compare algorithm from the lecture to algorithms their algorithms:

- (a) *partial_size1(current_length, level)*
- (b) *partial_size2(current_length, prev_char, level)*
- (c) *Is a table size of max_length · max_level enough?*

1. In the lecture l stands for level / rank, while in the note, it stands for the length of a password
2. In the lecture p_1 or probability(e_1) stands for the probability of an char / edge, while in the note they use ν (Ny)
3. For the probability of a path / password, they use the symbol θ (Theta)
4. They don't show how to save the save the information, which is goatherd with partial_size1. While algorithm 1 does show that the information gets stored in a table.
5. The partial_size1 algorithm supposes that the Markov model is of order 1. While the Algorithm 1 does not.
6. The partial_size1 and 2 algorithms only allow passwords with a fixed (same) length, while the algorithm from the lecture allows passwords of different sizes.
7. Also the table size of max_length · max_level is not enough. In the paper for partial_size1 they argue that a 2D-Array of size l times number of levels should be sufficient. That is the case for partial_size1, as each probability is independent of one another. But not for partial_size2, because it works for Markov models order 1, thus it loses the independence of probability and a larger table will be needed.